# Software work in mathematical biology

## Maxie Dion Schmidt

maxieds@gmail.com
http://people.math.gatech.edu/~mschmidt34
https://github.com/maxieds
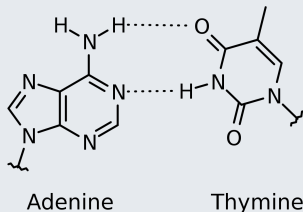
Sandia National Labs
Technical Presentation
Spring 2022

# Applications of RNA research

- ▶ Ribonucleic acid (RNA) sequencing has been utilized in many aspects of cancer research and therapy
- ▶ mRNA vaccines have elicited potent immunity against infectious disease targets in animal models of several key viruses in the past few years
- ▶ Research on RNA may eventually play a central role in diagnostics, therapeutics and research in medical applications
- ▶ As the cellular roles for RNA molecules continue to grow, so does the importance of gaining functional insight from structural analyses

# RNA sequence basics

▶ RNA is a single-stranded molecule similar to DNA; it carries messenger instructions to the double-stranded DNA that encodes genetic instructions required for life processes

▶ A strand of RNA has a backbone comprised of alternating sugar (ribose) and phosphate groups

▶ Each sugar has one of four base types attached to it:
Adenine – **A**, uracil – **U**, cytosine – **C**, or guanine – **G**

▶ Each of the **A–U–C–G** bases can fold to form bonds with one another

---

**Example:** *The adenine-thymine (AT) Watson-Crick base pair (*`WP/Base_pair`*):*



Adenine          Thymine

# RNA secondary structures

- ▶ A RNA for an organism is defined by the 1D base sequence and can have more than one secondary and tertiary structure by folding

- ▶ While obtaining the 1D sequence information is now relatively easy via sequencing, characterizing 3D molecular conformations is still comparatively hard

- ▶ Hence, understanding the 2D secondary structures, i.e. the intra-sequence base pairing, remain a crucial component of ribonomics research

- ▶ RNA folding prediction programs (e.g., GTFold or RNAFold) can take a base sequence and use probabilistic algorithms to generate the secondary structure

- ▶ The secondary structures generated by such computational procedures can be influenced by certain pairing penalties and by environmental constraints like a thermodynamic model specification

# What is mathematical biology?

▶ Mathematical biology (abbreviated MathBio) is a field that uses mathematical models and theoretical abstractions of the structure of living organisms

▶ These models are used to explore biological principles dictating the underlying structures, the development and the behavior of these living systems

▶ In this talk, we will discuss my work as a software engineering RA with the *Georgia Tech Discrete Mathematics and Molecular Biology* (gtDMMB) research group led by Christine Heitsch

▶ N.b., my official title within the group is endearingly designated as the "*Code Goddess*" :)

# The RNAStructViz application

# RNAStructViz: Graphical base pairing analysis

▶ The historical RNAStructViz application was a project developed by Professor Christine Heitsch and her husband (Ph.D. in computer graphics from UCB) to visualize RNA secondary structures

▶ By the time I arrived and started work with gtDMMB in the summer of 2018, the old C++ source code was badly broken with modern Linux and MacOS compilers

▶ My work was to modernize the C++ source, add support for enhanced graphics using the cairo library, and to generally improve and support the project in the long term

▶ The key feature RNAStructViz provides is visualization and comparisons with arc diagrams
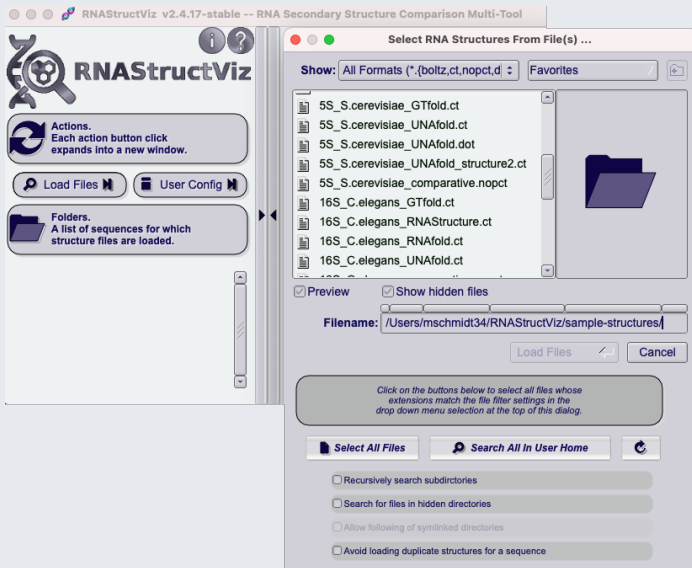
# RNAStructViz – Comparison of features

| ↓ Feature sets | RNAStructViz | FORNA | jViz.RNA | R-chie | RNAbows | VARNA |
|---|---|---|---|---|---|---|
| **Software support →** | | | | | | |
| **❶ — Platform and availability —** | | | | | | |
| Mac OSX support | ✔ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Linux / Unix support | ✔ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Windows support | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Open source software | ✔ | ✓ | ✓ | ✓ | ✗ | ✓ |
| Requires external libraries | ✔ | ✓ | ✓ | ✓ | ✗ | ✓ |
| **❷ — Software usability criteria —** | | | | | | |
| Graphical user interface | ✔ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Web interface | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ |
| Multi-window interface | ✔ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Compares 2 structures at once | ✔ | ✓ | ✗ | ✓ | ✓ | ✗ |
| Compares 3 structures at once | ✔ | ✓* | ✗ | ✗ | ✗ | ✗ |

| ↓ Feature sets | RNAStructViz | FORNA | jViz.RNA | R-chie | RNAbows | VARNA |
|---|---|---|---|---|---|---|
| **Software support →** | | | | | | |
| **❸ — Support for standard formats —** | | | | | | |
| CT files | ✔ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Dot-bracket files | ✔ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Built-in file viewer | ✔ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Requires specialized format | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ |
| Can edit sequence data | ✗ | ✓ | ✓ | ✗ | ✓ | ✓ |
| **❹ — Views and diagram type support —** | | | | | | |
| Has comparison statistics | ✔ | ✗ | ✓** | ✓ | ✓ | ✗ |
| Plots circular arc diagrams | ✔ | ✗ | ✓ | ✓ | ✓ | ✗ |
| Plots radial diagrams | ✔ | ✓ | ✓ | ✗ | ✗ | ✓ |

*A comparison of selected features across related tools; an extended survey appears in the RNAStructViz WIKI.*

# RNAStructViz Screenshot – Loading sample structures I

# RNAStructViz Screenshot – Loading sample structures II

# RNAStructViz Screenshot – Arc diagram window

# Arc diagram window – Discussion

▶ The bases indexed from position #1 to `#LengthOfBaseSequenceString` are placed at equidistant spacings around a circle

▶ The sequentially numbered base pairs are ordered around the circle counter-clockwise starting from the position at 270 degrees, or $\frac{3\pi}{2}$ radians, labeled by the 5′ | 3′ directional arrows in the display window

▶ An arc connecting paired bases is drawn within the circle

**Example:** *A furanose (sugar-ring) molecule with carbon atoms labeled using standard notation. The 5′ is upstream, whereas the 3′ is downstream – diagram taken from* `WP/Directionality_(molecular_biology)`*:*

# Arc diagrams – Discussion example – S. Cerevisiae (yeast)

GGUUGCGGCCAUAUCUACCAGAAAGCACCGUUUCCCGUCCGAUCAACUGUGUUAAGCUGGUAGA
((((((((((....(((((((((...((((.(((((......))..))).))))....))))))))))
GCCUGACCGAGUAGUGUAUGGGUGACCAUACGCGAAACUCAGGUGCUGCAAUCU
.....(((((((.(((((((((....)))))))))...))))).))))))))))).



**(Actual microscopic views)**     **(Radial view of 2D MFE structure)**

# Arc diagram window – Zoom select

# Arc diagram zoom – Radial layout visualization

# Arc diagram zoom – CT segment visualization

# Arc diagram window – Comparing multiple structures

# RNAStructViz Screenshot – Statistics window

# GTFoldPython software project

# GTFold – Overview and algorithms

▶ Accurate and efficient RNA secondary structure prediction remains an important open problem in computational molecular biology

▶ GTFold is the first implementation of RNA secondary structure prediction by thermodynamic optimization for modern multi-core computers

▶ The difference in the parallel prediction program used by GTFold to take full advantage of today's modern computing technology is particularly valuable to researchers working with very RNA sequences, such as RNA viral genomes

# GTFold – Overview and algorithms II

▶ The decade old historical source code for GTFold produced several command line only utilities that could be used to generate secondary structures (especially, MFE and MFE structures)

▶ We were motivated by the need to run Python3 on a webserver to index sequences and corresponding computational data generated with GTFold

▶ One of my ideas as a gtDMMB graduate research assistant in the fall of 2019 was to write Python3 bindings around the original GTFold sources in C++

▶ Perhaps the most impressive slide of source code in the project is it's highly robust and cross platform war-tested `Makefile` and supporting set of install scripts :)

# GTFoldPython – Introduction

► The GTFoldPython (GTFP) project provides Python3 bindings based around the historical GTFold sources in C++
► The backend uses the Python3 C API
► The frontend is a wrapper library that uses CTypes to call the C API functions from a dynamic shared library (e.g., Windows DLL / MacOS Dylib / Linux SO)

# GTFoldPython – Comparison – Python C API code

```
1   PyObject * GetMFEStructure(const char *baseSeq, ConsListCType_t consList, int consLength) {
2        /* Error checking omitted ... */
3        MFEStructRuntimeArgs_t rtArgs;
4        InitMFEStructRuntimeArgs(&rtArgs);
5        rtArgs.baseSeq = baseSeq;
6        SetRTArgsSequenceLength(rtArgs, strlen(baseSeq));
7        if(ParseGetMFEStructureArgs(consList, consLength, &rtArgs) != GTFPYTHON_ERRNO_OK) {
8            FreeMFEStructRuntimeArgs(&rtArgs);
9            return ReturnPythonNone();
10       }
11       if(InitGTFoldMFEStructureData(&rtArgs) != GTFPYTHON_ERRNO_OK) {
12           FreeMFEStructRuntimeArgs(&rtArgs);
13           return ReturnPythonNone();
14       }
15       double mfe = ComputeMFEStructure(&rtArgs);
16       if(GetLastErrorCode() != GTFPYTHON_ERRNO_OK) {
17           return ReturnPythonNone();
18       }
19       if(WRITEAUXFILES) {
20           ConfigureOutputFileSettings();
21           save_ct_file(outputFile, baseSeq, mfe);
22       }
23       char *dbMFEStruct = ComputeDOTStructureResult(rtArgs.numBases);
24       PyObject *mfeTupleRes = PrepareMFETupleResult(mfe, dbMFEStruct);
25       Free(dbMFEStruct);
26       FreeMFEStructRuntimeArgs(&rtArgs);
27       FreeGTFoldMFEStructureData(rtArgs.numBases);
28       if(mfeTupleRes == NULL) {
29           return ReturnPythonNone();
30       }
31       return mfeTupleRes;
```

# GTFoldPython – Comparison – Wrapper library code

```
1    ## Library initialization code:
2    if GTFPConfig.PLATFORM_DARWIN:
3        GTFoldPython._libGTFoldHandle = ctypes.cdll.LoadLibrary("GTFoldPython.dylib")
4    else:
5        GTFoldPython._libGTFoldHandle = ctypes.PyDLL("GTFoldPython.so",
6                                          mode=ctypes.RTLD_GLOBAL, use_errno=True)
7    @staticmethod
8    def _WrapCTypesFunction(funcname, restype=None, argtypes=None):
9        return GTFoldPython._libGTFoldHandle.__getattr__(funcname)
10
11    @staticmethod
12    def GetMFEStructure(baseSeq, consList = []):
13        """Get the MFE and MFE structure (in DOTBracket structure notation)
14        :param baseSeq: A string of valid bases (ATGU/X)
15        :param consList: A list of constraints on the MFE structure
16        :return: A tuple (MFE as double, MFE structure as string in DOTBracket notation)
17        :rtype: tuple
18        """
19        GTFoldPython._ConstructLibGTFold()
20        resType = ctypes.py_object
21        argTypes = [ GTFPTypes.CStringType,
22                     GTFPTypes.FPConstraintsListType(consList),
23                     ctypes.c_int ]
24        libGTFoldFunc = GTFoldPython._WrapCTypesFunction("GetMFEStructure", resType, argTypes)
25        (mfe, mfeStruct) = libGTFoldFunc(GTFPTypes.CString(baseSeq),
26                                         GTFPTypes.FPConstraintsList(consList),
27                                         len(consList))
28        return (float(mfe), str(mfeStruct))
```

# GTFoldPython – Example – Find MFE and MFE structure

**External Python3 script source:**

```
1   import sys, os
2   from GTFoldPythonImportAll import *
3
4   GTFP.Init()
5   GTFP.Config(quiet = False, debugging = False, verbose = False, stdmsgout = "stderr")
6
7   baseSeqFPCons = "GCAUUGGAGAUGGCAUUCCUCCAUUAACAAACCGCUGCGCCCGUAGCAGCUGAUGAUGCCUACAGA"
8   consListFP = GTFPUtils.ReadFPConstraintsFromFile("../Testing/TestData/tRNA/yeast.fa.cons")
9
10  (mfe, mfeDOTStruct) = GTFP.GetMFEStructure(baseSeqFPCons, consListFP)
11  print("MFE_%1.3f_=>_MFE_DOT_STRUCT_\"%s\"\n\n" % (mfe, mfeDOTStruct))
```

**Terminal output printed upon invoking the script above:**

```
1   MFE -17.200 => MFE DOT STRUCT "((((((((((.........))))............((((.......)))))....)))))......
```

# Thermodynamic models – Introduction I

▶ Standardized sets of empirically well tested nearest neighbor (NN) parameters can be used to predict the stability of RNA secondary structures (and hence liklihood of a given structural variant to exist naturally)

▶ Common sets of these parameters, such as the NNDB models – Turner1999 and Turner2004 – e.g., those documented in [7, 1] – exist and are publicly posted for download on the pioneering researchers' websites

▶ There are also programs for generating modified sets of these parameters to use in computational simulation for comparison against benchmarked data:
E.g., see https://github.com/gtDMMB/GTModify

# Thermodynamic models – Introduction II

▶ **Key idea:** Model RNA folding as a *dynamic system* of structures which represent the states of the system at a given time

▶ After a very long time, asymptotically, we expect that a given base sequence $\mathcal{S}$ will form every possible structure $\mathcal{P}$

▶ For each sequence $\mathcal{S}$ and possible structure $\mathcal{P}$, there is an associated probability for observing $\mathcal{S}$ folding to $\mathcal{P}$ at any given time

▶ The idea is to use models from mathematics and thermal physics to aim to define the most _reasonable_ probabilities: $\mathbb{P}[\mathcal{P} \mid \mathcal{S}]$

# Thermodynamic models – The Boltzmann distribution I

- Let $\mathcal{X} := \{X_1, X_2, \ldots, X_N\}$ denote a system of states where the $i^{th}$ state, $X_i$, has an associated energy $E_i$

- We say that the system $\mathcal{X}$ is *Boltzmann distributed with temperature T* if and only if

$$\mathbb{P}[X_i] = \exp\left(-\beta E_i\right) \cdot Z^{-1}, \text{ for } Z := \sum_{i \geq 1} \exp\left(-\beta E_i\right), \beta := \left(k_B T\right)^{-1}$$

- The empirical value of the constant $k_B \approx 1.38 \times 10^{-23}$ is called the (classical) *Boltzmann constant*

- The parameter $\beta$ is called the *inverse temperature* of the system

# Thermodynamic models – The Boltzmann distribution II

▶ The Boltzmann distribution is usually attained in limiting *thermodynamic equilibrium* that occurs after letting the system progress for a long time period

▶ Why the dependence on a fixed temperature $T$ above?
  1. At very high temperatures we assume that all states are approximately equally probable to occur
  2. Whereas at very low temperatures we assume that only the "best" states are likely to occur

▶ Why is the Boltzmann distribution important here? We can give a proof of the following information theoretic fact:
  *The Boltzmann distribution (BD) makes the least number of assumptions on the model (requires the least external / extraneous information). That is, formally, the BD is the unique distribution with the lowest information content, or equivalently, with the maximal (Shannon) entropy.*

# Thermodynamic models – NN model – Definitions

▶ We define the *probability of a base pair* $(i, j)$ *for* $\mathcal{S}$ to be

$$\mathbb{P}[(i, j) \mid \mathcal{S}] := \sum_{(i,j) \in P} \mathbb{P}[P \mid \mathcal{S}]$$

▶ We can then derive (see [3]) a recursive NN prediction for individual base pairs in the optimal structures as follows:
  - Initialize: $N_{i,j} := 0$ for all $i - j \leq m$;
  - Recursively compute:

$$N_{i,j} = \max\left( N_{i,j-1}, \max_{i \leq k < j - m} N_{i,k-1} + N_{k+1,j-1} + 1 \right)$$

▶ Note that here, the notation $N_{i,j}$ denotes the distance from base pair $(i, j)$ (at the $i^{th}$ position) to the next base it should be paired with

# Thermodynamic models – NN and dynamic programming

▶ We define the *partition function $Z_{\mathcal{P}}$ for the set of RNA structures $\mathcal{P}$ of $\mathcal{S}$* by

$$Z_{\mathcal{P}} := \sum_{P \in \mathcal{P}} \exp(-\beta E(P))$$

▶ Computational strategy: Compute the partition functions $Z_{\mathcal{P}_{ij}}$ recursively

▶ This computation can be made efficient through the use of dynamic programming

# GTFoldPython – Specifying custom thermodynamic models

▶ This functionality used to be available in GTFold by setting:
  export GTFOLD_DATADIR=$(readlink -f /path/to/thermo/data)

**External Python3 script source:**

```
1   import sys, os
2   from GTFoldPythonImportAll import *
3
4   GTFP.Init()
5   GTFP.Config(quiet = False, debugging = False, verbose = False, stdmsgout = "stderr")
6
7   ## Configure the energy model and thermodynamic parameters
8   ## NOTE: Paths relative to the current running script location
9   #GTFP.SetGTFoldDataDirectory("../Testing/ExtraGTFoldThermoData/GTFoldTurner04/")
10  #GTFP.SetThermodynamicParameters(TURNER04, None)
11  GTFP.SetThermodynamicParametersFromDefaults(TURNER04)
12
13  baseSeqFPCons = "GCAUUGGAGAUGGCAUUCCUCCAUUAACAAACCGCUGCGCCCGUAGCAGCUGAUGAUGCCUACAGA"
14  consListFP = GTFPUtils.ReadFPConstraintsFromFile("../Testing/TestData/tRNA/yeast.fa.cons")
15
16  (mfe, mfeDOTStruct) = GTFP.GetMFEStructure(baseSeqFPCons, consListFP)
17  print("MFE_%1.3f_=>_MFE_DOT_STRUCT_\"%s\"\n\n" % (mfe, mfeDOTStruct))
```

**Terminal output printed upon invoking the script above:**

```
1   MFE -14.770 => MFE DOT STRUCT "((((((((........))))...........(((((......)))))....)))))......
```

# Concluding remarks

# The End

Questions?

Comments?

Feedback?

# Thank you for your time!

# References I

M. Andronescu, A. Condon, D. H. Turner and D. H. Matthews. *The determination of RNA folding nearest neighbor parameters*. Methods Mol. Biol. **1097**, pp. 45–70 (2014). `https://doi.org/10.1007/978-1-62703-709-9_3`

*GTFoldPython software documentation*. `https://github.com/gtDMMB/GTFoldPython/wiki`

MIT Math 18.417 Course Lecture Notes. *RNA ensembles*. Available online: `https://math.mit.edu/classes/18.417/Slides/rna-ensembles.pdf`

*RNAStructViz software documentation*. `https://github.com/gtDMMB/RNAStructViz/wiki`

Schmidt, M. D., Kirkpatrick, A., and Heitch, C. *RNAStructViz: graphical base pairing analysis*. Bioinformatics **197** (2021). `https://doi.org/10.1101/2021.01.20.427505`

# References II

📄 Swenson, M. S., Anderson, J., Ash, A., Gaurav, P., Sukos, Z., Bader, D. A., Harvey, S. C., and Heitsch, C. E. *GTfold: Enabling parallel RNA secondary structure prediction on multi-core desktops*. BMC Research Notes. 5(1): 341, 2012. https://github.com/gtDMMB/gtfold

📄 D. H. Turner and D. H. Matthews. *NNDB: the nearest neighbor parameter database for predicting stability of nucleic acid secondary structure*. Nucleic Acids Research, Volume 38, D280–D282 (2010). https://doi.org/10.1093/nar/gkp892